

ESnet and ANI Testbed Update

Brian Tierney

DOE PI Meeting

Bethesda, MD

March 2, 2011





Outline



- 1. ESnet Overview
- 2. ESnet R&D projects: OSCARS and perfSONAR
- 3. ANI Testbeds

ESnet is a Mission Organization





DOE Office of Science mission: to deliver scientific discoveries and major scientific tools for transforming our understanding of nature, and to advance the energy security, economic security, and national security of the United States.

ESnet mission: to accelerate scientific discovery for the DOE Office of Science by delivering unparalleled network infrastructure, services, tools, and innovation.

ESnet Supports DOE Office of Science



SC provides broad support for Labs, Facilities, people

- almost \$5B/year in funding
- 45% of Federal support for physical sciences research
- key funding for basic research in biology, computing, energy, climate
- over 100 Nobel Prizes in past 60 years

Supporting > 27,000 PhDs, grad students, engineers at >300 institutions.

Provides world's largest collection (32) of scientific user facilities:

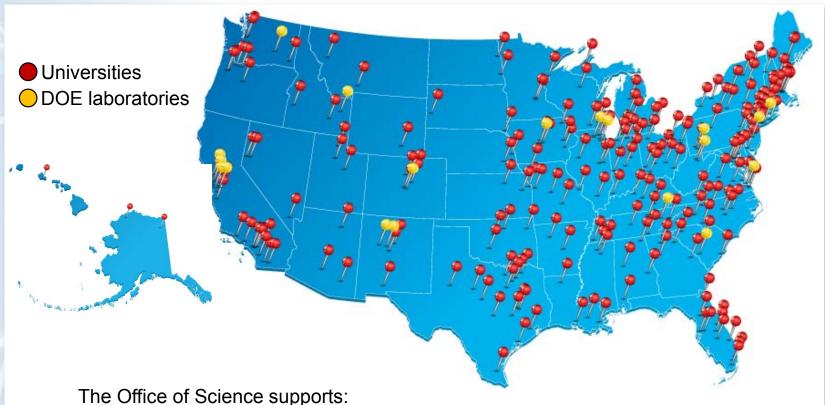
 supercomputer centers, accelerators, light sources, neutron sources, electron microscopes, nano-scale centers, a sequencing center, fusion facilities.

ESnet is one of them – connecting DOE sites, facilities, scientists and collaborators.

- optimized for science data transport
- every service exists to support scientific discovery

ESnet Supports DOE Office of Science





- 27,000 Ph.D.s, graduate students, undergraduates, engineers, and technicians
- 26,000 users of open-access facilities
- 300 leading academic institutions
- 17 DOE laboratories

ESnet is Embedded in a Global Research Networking Ecosystem





Lambda Integrated Facility Visualization by Robert Patterson, NCSA, University of Illinois at Urbana-Champaign Data Compilation by Maxine D. Brown, University of Illinois at Chicago Texture Retouch by Jeff Carpenter, NCSA Earth Texture, visibleearth.nasa.gov www.glif.

ESnet is Embedded in a Global Research Networking Ecosystem



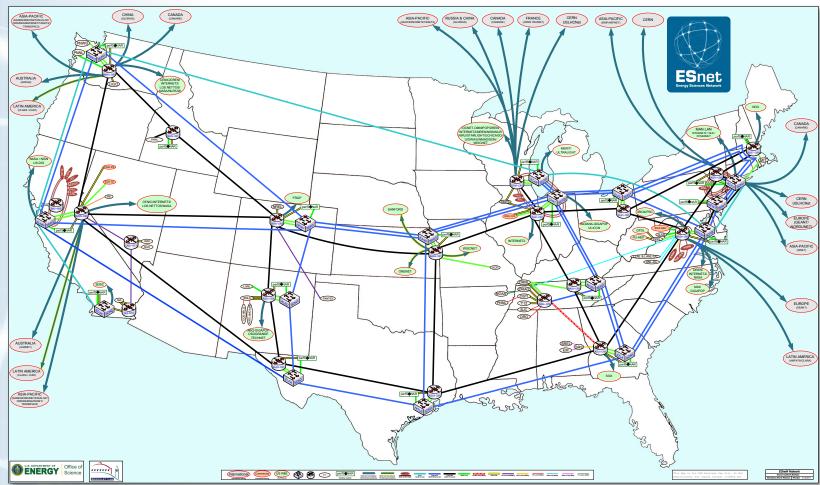
Most important implication: requirement for *multi-domain* coordination.

- Large-scale science is a multi-domain, end-to-end endeavor.
- We must build services across networks, exchanges, campuses.
 - Problem resolution, data mobility, measurement / monitoring, guaranteed services, collaboration – all multi-domain!
- We spend much time coordinating, communicating
 - International: GLIF, DICE, LHCOPN, LHCONE, TERENA
 - US: Joint Techs, ESCC, Supercomputing, many more

Our drivers are different from those of commercial providers.

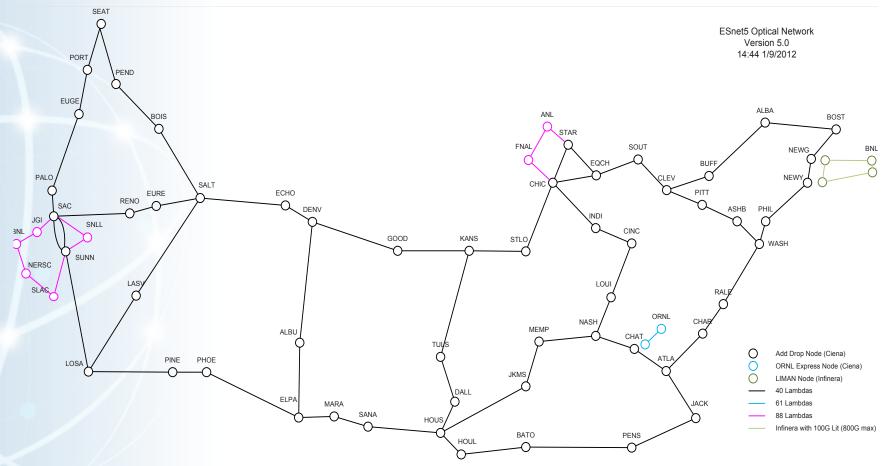
ESnet4 Topology (n x 10G core)





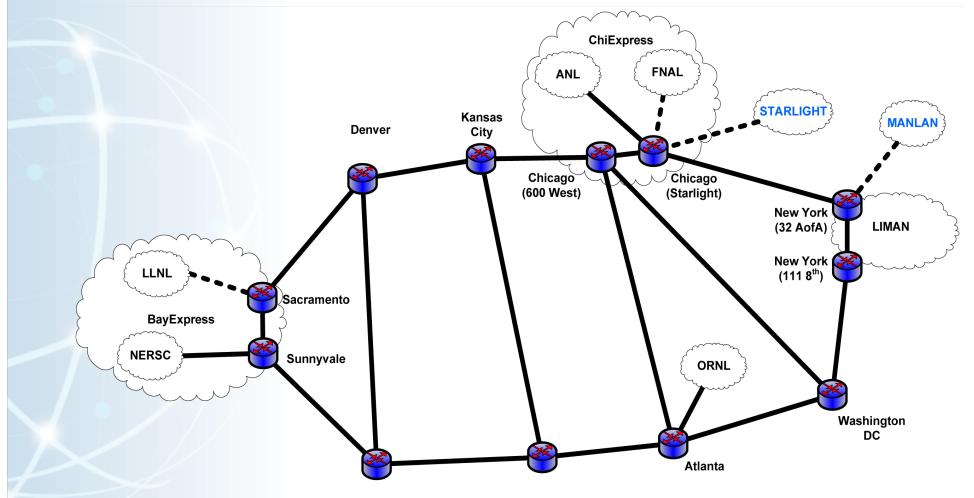
ESnet5 Optical Infrastructure Footprint





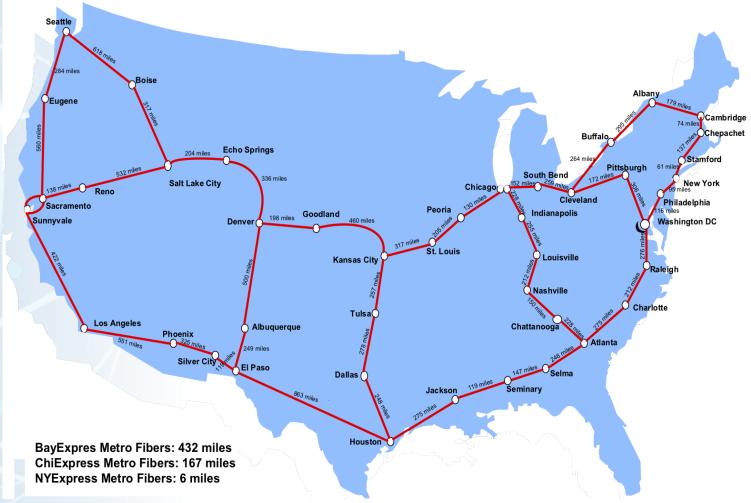
ESnet5 L2/3 Architecture (n x 100G core)





12,924 miles of Long-Haul Dark Fiber, ESnet5+

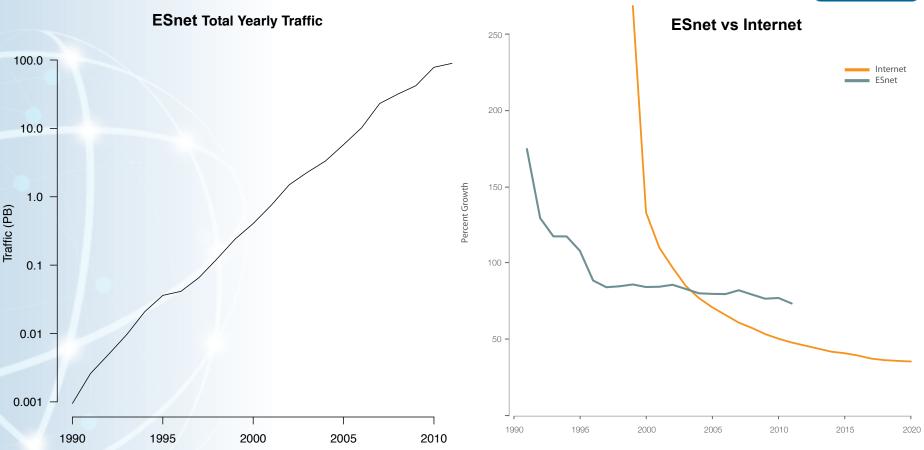




One Important Driver for Innovation



12



PB Accepted per Year

3/2/12

Compound Annual Growth Rates for ESnet [measured] and Internet

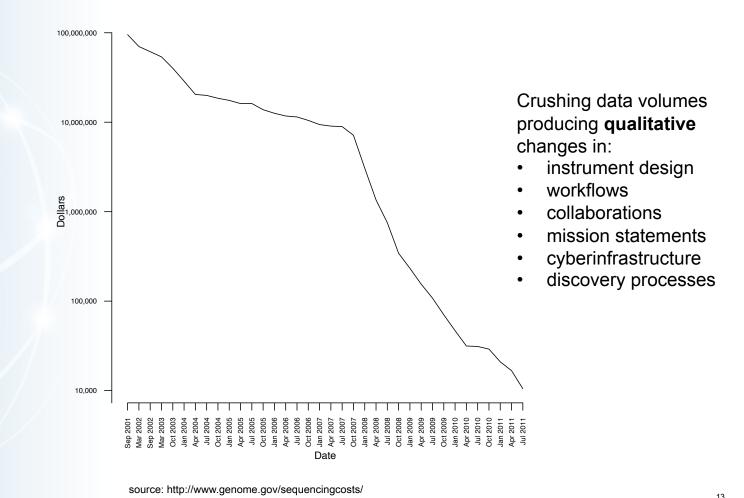
core [estimated and predicted, see "Power Trends in Communication Networks", IEEE Journal of Selected Topics in Quantum Electronics,

VOL. 17, NO. 2].

ESnet is Engineered for the Data Explosion Flood Tsunami Revolution



Cost of Sequencing per Genome



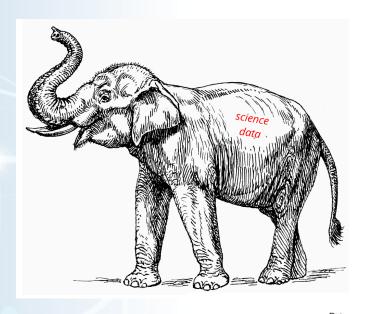
Lawrence Berkeley National Laboratory

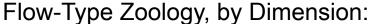
3/2/12

U.S. Department of Energy | Office of Science

Elephant (or Alpha) Flows







- size: elephant and mouse
- rate: cheetah and snail
- duration: tortoise and dragonfly
- burstiness: porcupine and stingray

Kun-chan Lan and John Heidemann, "A Measurement Study of Correlations of Internet Flow Characteristics." *ACM Comput. Netw.* 50, 1 (January 2006), 46-6.

Or:

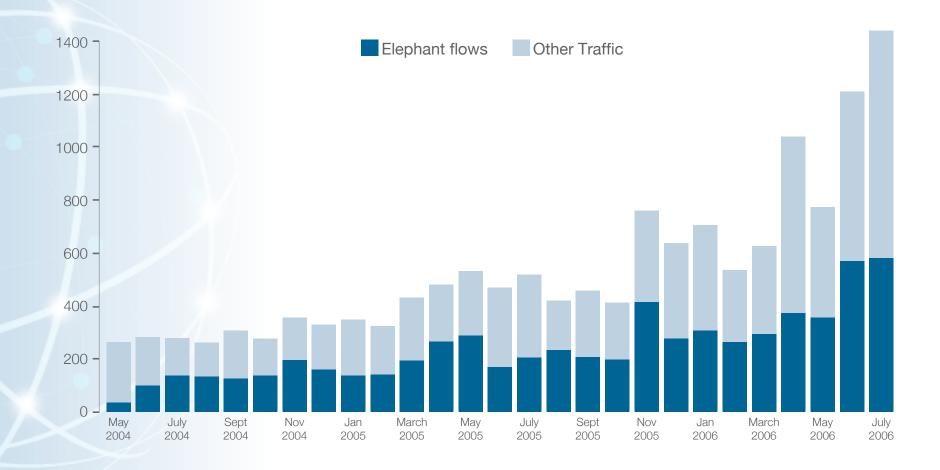
- Alpha flows (> H byes in < T sec)
- caused by transfers of large files over fast links

S. Sarvotham, R. Riedi, and R. Baraniuk, "Connection-Level Analysis and Modeling of Network Traffic," *ACM SIGCOMM Internet Measurement Workshop*, November 2001, 99–104.

Science Traffic is Different

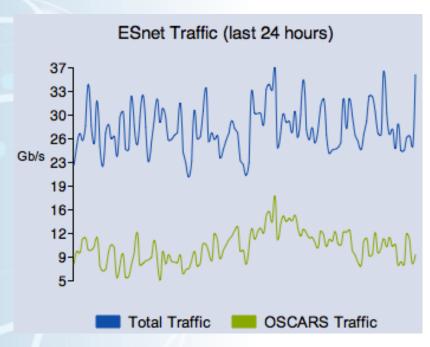


15

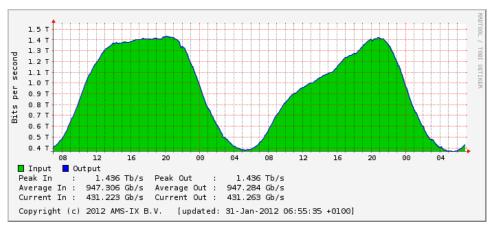


Science Traffic is Different





ESnet Accepted Traffic on January 28, 2012 (one day) http://www.es.net

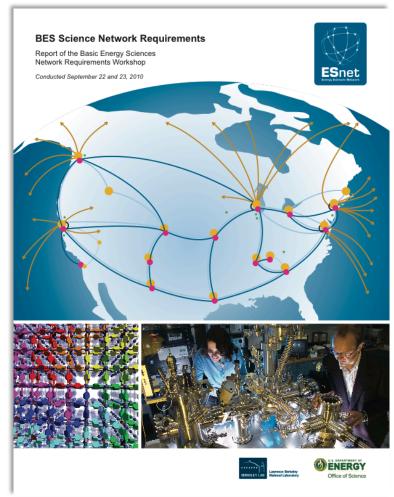


http://www.ams-ix.net/statistics/

How do we Know what our **Customers Need?**



- Each Program Office at Office of Science has a dedicated requirements workshop devoted to its needs
- Two requirements workshops per year, with attendees chosen by science programs
- Discussion centered on science case studies
- Requirements derived from case studies + discussions
- Reports contain requirements, case study text





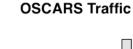
ESnet R&D

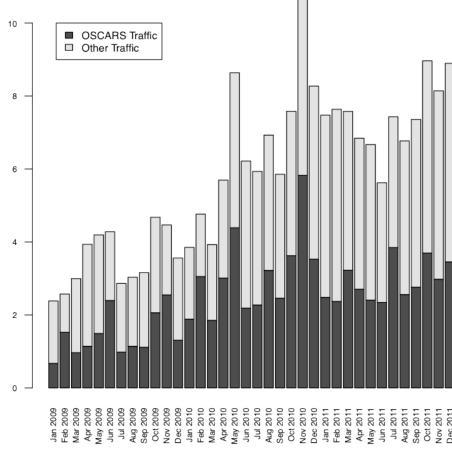
ESnet Research Success: OSCARS



OSCARS: an application and protocol suite for dynamic, interdomain circuit management

- notable ESnet innovation
- used for service guarantees, traffic engineering, overlays
- carries almost ~50% of ESnet traffic (alpha flows)
- a customer-facing service: web interface for people, API for apps
- deployed in 21+ networks





Production services OSCARS



Research Motivation

Network guarantees needed by service sensitive applications & large flows

Objective

Provide secure end-to-end guaranteed circuits across ESnet & multiple domains.

2003: Flow analysis and identification of large flows

2007: OSCARS production service 2008: GLIF interops multiple implementations using FENIUS (ESnet driven) 2011: Standards based NSI protocol defined in OGF, successful interop with multiple independent implementations



2008: DICE collaboration completes IDC protocol v1.0

2009: ARCHStone research project funded

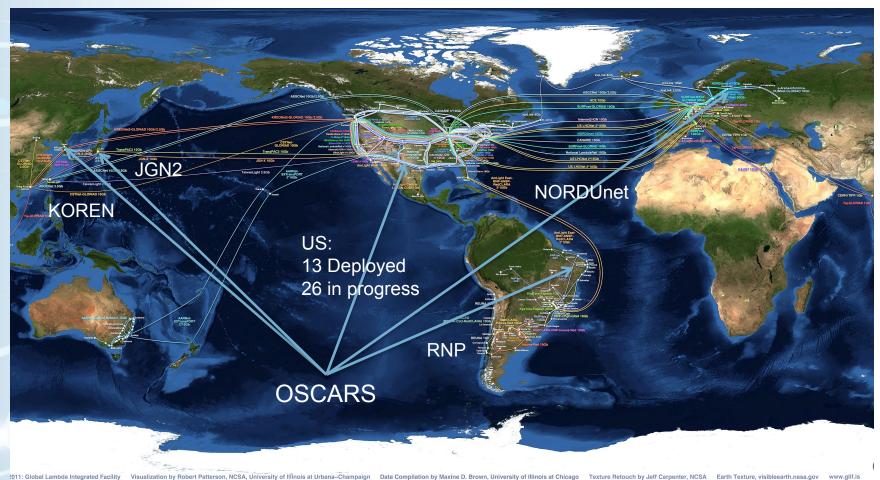
Impact

- Carries 50% of ESnet's 100+PB yearly data
- Deployed in about 21 networks
- Under evaluation by an additional 26 greenfield deployments

OSCARS Deployments Worldwide

WAN, Regional, Campus, Testbeds





ESnet Research Success: perfSONAR



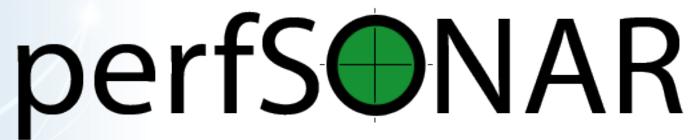
Distributed platform and protocol suite for performance measurement, monitoring, and visualization.

Assures end-to-end health of global networks.

400+ nodes deployed worldwide.

Developed in a global collaboration of R&E networking organizations.

Helps us find and fix problems!



Network tools perfSONAR



Research Motivation

 Certain performance problems only apparent to long-latency TCP flows, testing those requires source/sink in other domains

Objective

 Provide an infrastructure that identifies and isolates performance problems in a multi-domain environment.

2000: Formation of OGF group

2005: nonstandard GE performance testers

2007:perfSONAR schema published

2009-11: Deployment at all 10G ESnet sites. perfSONAR release every 6 months

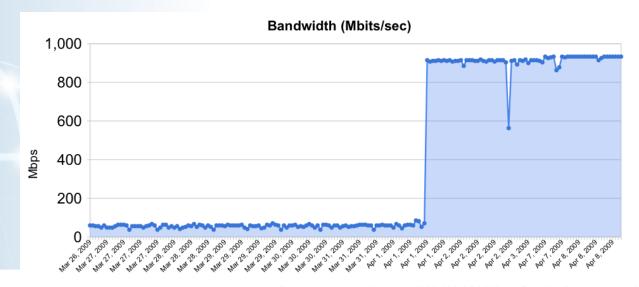
2003: First OGF doc on network measurements 2006: 1st release of software. perfSONAR system 10GE b/ w and GE OWAMP 2008: 1st release of perfSONAR-PS. 20 nodes deployed

Impact

- Almost every ESnet site has fixed network performance problems
- LHC collaboration has deployed extensively
- 400+ perfSONAR hosts deployed on 100 networks and in 14 countries

perfSONAR in Use

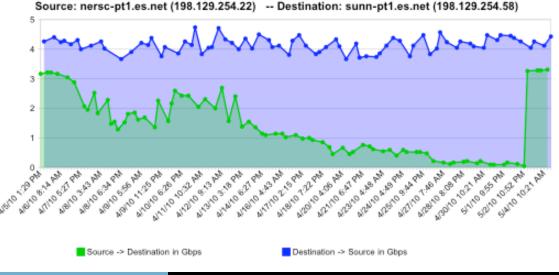




Effect of full routing table, followed by reboot.

Effect of gradual optical line card failure.

1/29/2012

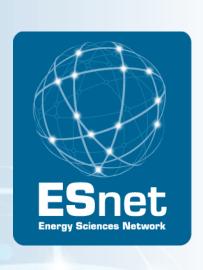


Creating VO-Specific Dashboards from perfSONAR Data



Status of perfSONAR Throughput Matrix

_	0	1	2	3	4	5	6	7	8
		-		3	7	3	•	'	
0:atlas-npt2.bu.edu	-	OK OK	OK OK	OK OK	OK OK	OK OK	UNKNOWN OK	OK OK	OK OK
		OK .	OK .	OK .	OK .	OK	OK .	OK .	OK .
1:lhcmon.bnl.gov	OK	-	OK	OK	OK	OK	OK OK	ОК	OK
	ОК		ОК	ОК	ОК	ОК	UK	UNKNOWN	ОК
2:ps2.ochep.ou.edu	ОК	<u>ок</u>	-	ОК	ОК	ОК	ОК	ОК	ОК
	ОК	ок		ОК	ОК	ОК	UNKNOWN	ОК	ОК
3:psmsu02.aglt2.org	ОК	ОК	ОК	-	ОК	ОК	UNKNOWN	ОК	ОК
	ОК	ОК	ОК		ОК	ОК	UNKNOWN	ОК	ОК
4:netmon2.atlas-	ОК	UNKNOWN	ОК	ОК	-	ОК	ОК	ОК	ОК
swt2.org	UNKNOWN	ОК	ОК	ОК		UNKNOWN	UNKNOWN	ОК	ОК
5:iut2-net2.iu.edu	ОК	ОК	ОК	ОК	ОК	-	ОК	ОК	ОК
	ОК	ОК	ОК	ОК	ОК		ОК	ОК	ОК
6:psnr-	ОК	ОК	UNKNOWN	UNKNOWN	UNKNOWN	ОК	-	ОК	UNKNOWN
bw01.slac.stanford.edu	UNKNOWN	ОК	ОК	UNKNOWN	UNKNOWN	ОК		ОК	UNKNOWN
7:uct2-	ОК	ОК	ОК	ОК	ОК	ОК	ОК	-	ОК
net2.uchicago.edu	ОК	ОК	ОК	ОК	ОК	ОК	ОК		ОК
8:psum02.aglt2.org	ОК	ОК	ОК	ОК	ОК	ОК	UNKNOWN	ОК	-
	ОК	ОК	ОК	ОК	ОК	ОК	UNKNOWN	ОК	



DOE's Advanced Network Initiative (ANI) Testbed Project Update

Testbed Value to Network Researchers



A realistic national-scale environment for innovative network research that would be impossible for most researchers to create in their labs

Many types of network research require realistic high-latency environments

Maximum flexibility: researchers get "super-user" access to everything

- ability to install custom OS on hosts
- ability to do full routers/switch configuration
- ability to install custom router software

A controlled environment that supports reproducible results

Priority is given to research that aligns with ESnet strategic focus areas, such as:

High performance network protocols, data transfer middleware, Openflow + MPLS integration strategies, energy efficient networking

ESnet Research Testbeds

Fertile environment for researchers to test and prove concepts

Control Plane Testbed

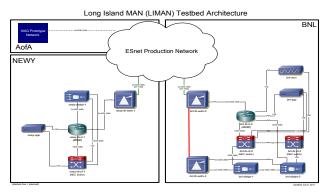
- Long Island fiber, Infinera transport
- Juniper, NEC OpenFlow, Data transfer nodes, VMs

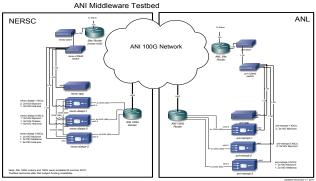
100G Testbed

- High-speed protocol research
- Routed network from Argonne to NERSC

Dark Fiber Testbed

- Disruptive network approaches
- Quantum-key encryption





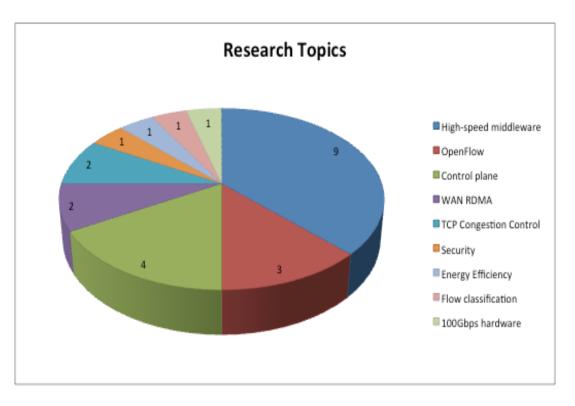


Community impact so far:

Peer-reviewed Testbed projects



- 25 projects accepted:
 - 17: testbed proposal review
 - 8: direct DOE/ASCR funding
 - 5 from Industry; 8 from DOE
 labs; 3 from NASA; 1 from
 DISA; 7 from Universities
- Strong diversity in research topics
- Four papers accepted or submitted for publishing, more to follow



Control Plane Testbed



Control Plane Testbed

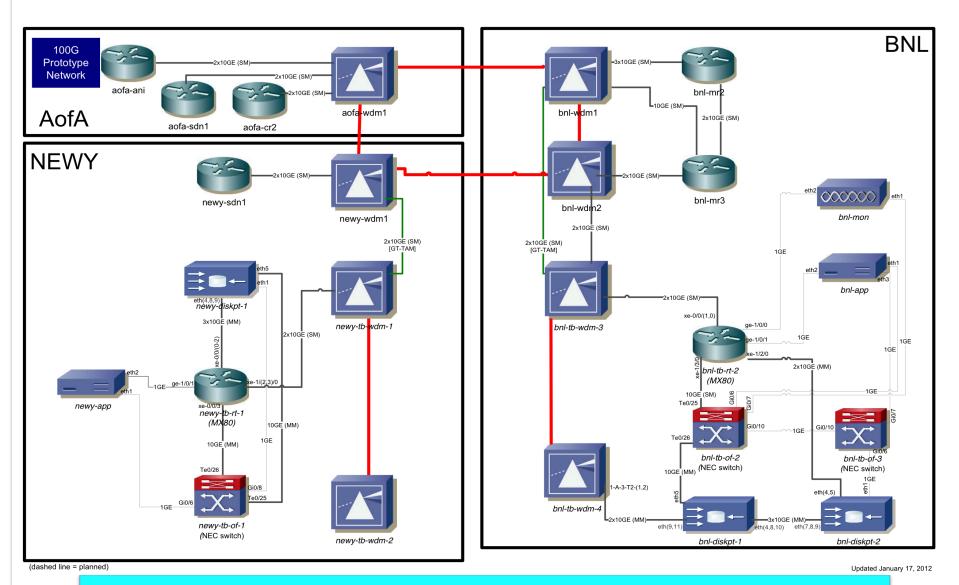
700 KM loop on Long Island Between Manhattan and Brookhaven National Lab

Infinera DWDM switches, NEC OpenFlow switches, Juniper router (soon with OpenFlow support), 6 hosts

Research Examples

- OpenFlow experiments
- Path Computation algorithms
- MPLS/OpenFlow Integration

Control Plane Testbed, Long Island, NY



Notes:

3/2/12 "App Host": can be used for researcher application, control plane control software, etc. Can support up to 8 simultaneous VMs

- -"I/O Testers" are capable of 15 G disk-to-disk or 35G memory-to-memory
- -Other infrastructure not shown: VPN Server, file server (NFS, webdav, svn, etc.)

100G Testbed



100G Testbed

100Gbps path from Oakland, CA (NERSC) to Argonne, IL (ANL)

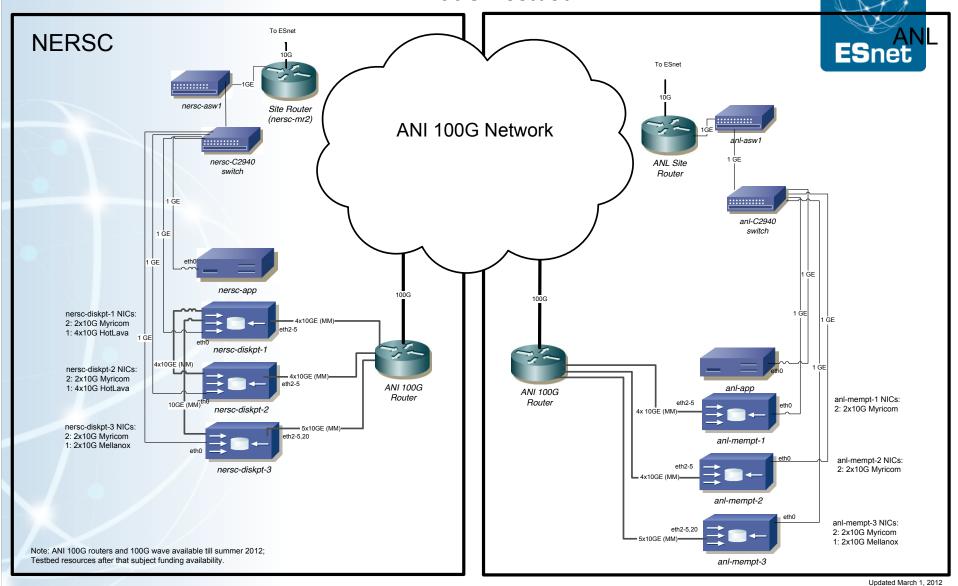
3 high-end hosts at each site, capable of >39 Gbps each

Research Examples

- 100G data transfer tools
 - FTP100, GridFTP, BestMan, Saratoga
- TCP alternatives
 - RDMA, RoCE
- TCP enhancements
 - PERT, RAPID TCP

Available as of Jan 3, 2012

ANI 100G Testbed



Related Resources



Magellan at NERSC

currently 15 nodes with 1 10GE, GPFS backend

soon: 8 nodes with 2x10GE, GPFS backend

Magellan at ANL

went away in January

ORNL test cluster

12 nodes, each with flash-based scratch disk

FNAL test cluster

in planning stage

Dark Fiber Testbed



Dark Fiber Testbed

Nationwide footprint (13,000 miles)

Researchers must pay to light the fiber, but no charge for access if results will be public

Research Examples

- Greentouch Energy experiments
- JPL quantum networking / quantum cryptography experiments

Dark Fiber Network







Current ANI Testbed Research



25 projects have been accepted to the testbed so far:

8 via direct DOE/ASCR funding

17 via testbed proposal review process

5 from Industry; 8 from DOE labs; 3 from NASA, 1 from DISA, 7 from Universities

Wide range of projects:

9: high-speed middleware 1: quantum communications

3: OpenFlow 2: Wide Area RDMA

4: other network control plane 2: TCP congestion control

1: 100Gbps end host hardware 1: security

1: network flow classification 1: energy efficiency

http://www.es.net/RandD/advanced-networking-initiative/current-testbed-research/

Potential Breakthroughs from Current Testbed Research



Unconditionally secure continuous variable quantum communications

Ability to scale TCP well beyond 10Gbps

Alternative transport protocols that scale better than TCP (e.g.: RDMA)

Ability to easily create end-to-end circuits

Ability to scale data transfer tools/middleware to 100Gbps and beyond 100Gbps host interface

For more information see:

http://www.es.net/RandD/advanced-networking-initiative/current-testbed-research/

Results to Date: Publications Submitted



- After only 7 months of being operational, there were already several publications submitted based on ESnet's Testbed results:
- Energy Profiling of Network Elements for Rate Adaptation Technologies, M. Ricca, A. Francini, S. Fortune, T. Klein, submitted to the ACM/IEEE 3rd International Conference on Future Energy Systems (e-Energy 2012).
- Traffic engineering in hybrid IP/optical circuit networks. Zhenzhen Yan, Chris Tracy, and Malathi Veeraraghavan, submitted to the IEEE 13th Conference on High Performance Switching and Routing (HPSR)
- Middleware Support for RDMA-based Data Transfer in Cloud Computing, Yufei Ren, Tan Li, Dantong Yu, Shudong Jin, and Thomas Robertazzi, submitted to IEEE International Parallel & Distributed Processing Symposium (IPDPS)
- Identifying gaps in Grid middleware on fast networks with the Advanced Network Initiative, G. Garzolglio and Dykstra, accepted at International Conference on Computing in High Energy and Nuclear Physics (CHEP 2012)

Results to Date: Press Releases



Joint Press Release by Orange Silicon Valley, Bay Microsystems and ESnet on the first wide area 40G RDMA results (http://goo.gl/57gSJ)

- Showed that RDMA moves data at up to 96 percent of the peak capacity of the network
 - Much more efficient than TCP
- Replicated at SC11 over a 10,000KM path (http://goo.gl/ZE0w8)

Greentouch press release on plans to use ESnet dark fiber testbed for energy experiments: http://goo.gl/OB6nD

Results to Date: Demonstrations



- ESnet Demonstration of OSCARS/OpenFlow integration on the ANI testbed at the Open Networking Summit held at Stanford.
 - http://goo.gl/4FFBg
- ESnet Demonstration of 9.9Gbps RDMA transfers from the ANI Testbed at BNL to Seattle at SC11

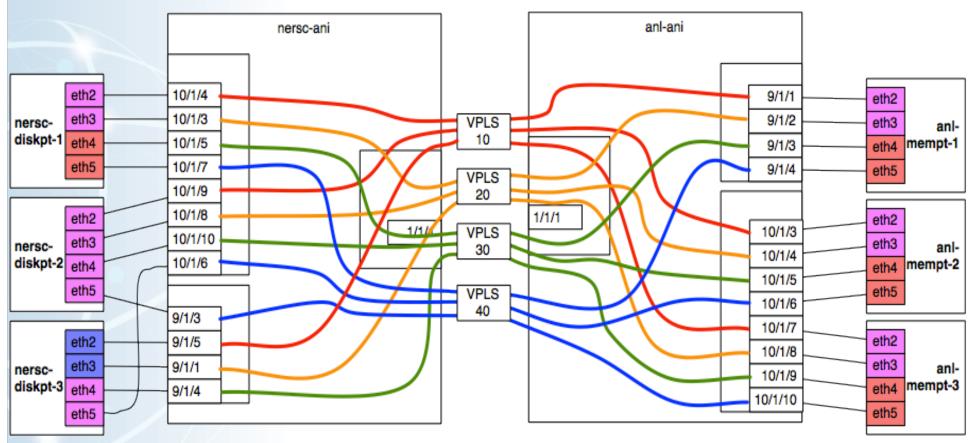
Results to Date: Industry Involvement



- Loaner hardware for 40GE testing
 - Infinera, Alcatel-Lucent (ALU), Bay Microsystems, Mellanox

100G Testbed Performance





48.6ms RTT, 97.9Gbps aggregate TCP throughput* with 10 TCP streams

3/2/12

43

Per-Stream Results



nersc-diskpt-1-v4012: nersc-diskpt-1-v4013: nersc-diskpt-1-v4014: nersc-diskpt-1-v4015: nersc-diskpt-2-v4012: nersc-diskpt-2-v4013: nersc-diskpt-2-v4014: nersc-diskpt-2-v4015: nersc-diskpt-3-v4014:	1179.1875 MB / 1179.2500 MB / 1179.1875 MB / 1179.1250 MB / 1179.2500 MB / 1179.0625 MB / 1179.3750 MB / 1179.1250 MB / 1179.1250 MB /	1.00 sec = 1.00 sec = 1.00 sec = 1.00 sec = 1.00 sec = 1.00 sec =	9891.8010 Mbps 9888.4787 Mbps 9891.1482 Mbps 9891.1581 Mbps 9891.9494 Mbps 9891.1580 Mbps 9893.1365 Mbps 9891.0690 Mbps	0 retrans 0 retrans 0 retrans 0 retrans 0 retrans 0 retrans 0 retrans
nersc-diskpt-3-v4015: nersc-diskpt-3-v4015: Octets Packets Errors Utilization (% of port	1121.8750 MB /	1.00 sec =	9410.9602 Nops 9410.9884 Mbps Input 162079 184615 0	0 retrans 0 retrans Output

Testbed Usage: Fully booked 24x7!

Tue

7 NERSC-ANL link, template 7 NERSC-ANL link, template 5 ANL Site Maintenance

Today < > February 2012

HNTES reserve newy-diskpt-1, bnl-diskpts, routers testing dark fibers and insert SLAC node 12 fiber testing from NERSC

Sun



Sat

6:30 ANL to NERSC testing

Fri

9 Jason's Networking Testir 6 NERSC-ANL link, GRC

12 fiber testing from NERSC 12 fiber testing from NERSC 12 fiber testing from NERSC 12 fiber testing from NERSC

Feb 2012

11 ANL-NERSC link reserve 1p ANL <-> NERSC 100G T 9 all hosts ANI+NERSC Terr 9p clim100 12p Engin's ANL-NERSC te 5p (No title) 10 11) 6 NERSC-ANL all hosts Ten HNTES reserve newy-diskpt-1, bnl-diskpt-2, routers 6 NERSC-ANL link, GRC 6 NERSC-ANL link, GRC 6 NERSC-ANL link, GRC 12:30 NERSC-ANL all hosts 12:30 Suny images Templa 6 all hosts template#1 6 nersc-diskpt-3, anl-diskpt- 7 nersc-diskpt-1, nersc-disk 5p Host testing 10 ANL-NERSC link: disk p 6 NERSC-ANL link, GRC 3:30p NERSC-ANL 12p anl-mempt-3, nersc-dis 12p ANL and NERSC hosts 10p Suny images Template 15 HNTES reserve newy-disk HNTES reserve newy-diskpt-1, bnl-diskpt-2, routers 12 Suny images Template # 6 ANL & NERSC all hosts w 6 ANL & NERSC all hosts w 6 ANL & NERSC all hosts w 6 all hosts w 6 all hosts w 12 all ANL&NERSC hosts, 1 9p climate100 all hosts tem 12p ANL & NERSC all hosts 8 fnal.anl-mempt-1.a, fnal.a 8 ANL-NERSC all nodes ten 6p ANL & NERSC all hosts vit 5p ANL & NERSC all hosts 11:30p nersc to ORNL only 3p clim100 all hosts templat 3p All hosts template#1 9p Sunv image all hosts ten 19 25 HNTES reserve newy-disk HNTES reserve newy-diskpt-1, bnl-diskpt-2, routers 8 Maintenance & Testing Fred Smith ANGEL reserve newy-tb-of-1 NEC switch 4p ANL/NERSC all hosts de 10G circuit AofA to NERSC, Hosts; bnl-diskpt-1, nersc-diskpt-1, nersc-diskpt-3 HNTES reserve newy-diskpt-1, bnl-diskpts, routers 6:30p climate100 all hosts t Presidents Day 4 10G circuit AofA to NERS(4 10G circuit AofA to NERS(12 NERSC/NAL all hosts im 6 FNAL, all performance ho 10 ANL/NERSC all hosts de 4 10G circuit AofA to NERS 6 NERSC-ANL link, GRC 6 nersc-diskot-2 all ANL hos 3 ANL-NERSC all nodes ten 4p ANL/NERSC all hosts de 5p All hosts reserved for ma 6 NERSC-ANL link, GRC 3p nersc-diskpt-2 anl-memp 4p ANL/NERSC all hosts de 11p all hosts template1 11:30p nersc & ANL all hos 11:30p clim100 image all ho 11p clim100 all hosts Mar 1 HNTES reserve newy-disk 12 fiber testing from NERSC 12 fiber testing from NERSC March 2012 Day Week Month 4 Days Agenda Mon Wed Fri Sat Sun Tue Thu 27 28 29 Mar 1 3

9 ANL 100G link maintenan 6 FNAL image on all perform

Feb 1

March 2012

HNTES reserve newy-dist 12 fiber testing from NERSC 12 fiber testing from NERSC ANGEL reserve bnl-tb-of-2

Testbed Access



Proposal process to gain access described at:

https://sites.google.com/a/lbl.gov/ani-testbed/

Testbed is available to anyone:

- DOE researchers
- Other government agencies
- Industry

Must submit a short proposal to the testbed review committee

Committee is made up of members from the R&E community and industry

Goal is to accept roughly five proposals every 6 month review cycle

Next round of proposals is due April 1, 2012

Future Testbed Challenges



Long term operational support of the Testbed

- Current testbed funding will run out in August of this year
 - Currently funded at 2FTE for testbed support
 - This is essential for continued success

Hope to secure funding to ensure persistent/indefinite availability

 100G wave from NERSC to ANL will transition to production use starting this summer, so need to identify funds to light a 2nd wave for testbed use

Dark Fiber will be available to the testbed until it is needed for production use

For some segments this may not be for several years.

More Information



http://sites.google.com/a/lbl.gov/ani-testbed/

email: ani-testbed-proposal@es.net

BLTierney@es.net